

How should a spatial-coverage sample design for a geostatistical soil survey be supplemented to support estimation of spatial covariance parameters?

R.M. Lark¹ *, B.P. Marchant

British Geological Survey, Keyworth, Nottinghamshire NG12 5GG, U.K.

1 **Abstract**

2 We use an expression for the error variance of geostatistical predictions, which in-
3 cludes the effect of uncertainty in the spatial covariance parameters, to examine the per-
4 formance of sample designs in which a proportion of the total number of observations
5 are distributed according to a spatial coverage design, and the remaining observations are
6 added at supplementary close locations. This expression has been used in previous studies
7 on numerical optimization of spatial sampling, the objective of this study was to use it to
8 discover simple rules of thumb for practical geostatistical sampling. Results for a range
9 of sample sizes and contrasting properties of the underlying random variables show that
10 there is an improvement on adding just a few sample points and close pairs, and a rather
11 slower increase in the prediction error variance as the proportion of sample points allo-
12 cated in this way is increased above 10 to 20% of the total sample size. One may therefore
13 propose a rule of thumb that, for a fixed sample size, 90% of sample sites are distributed
14 according to a spatial coverage design, and 10% are then added at short distances from
15 sites in the larger subset to support estimation of spatial covariance parameters.

16 **Keywords.** Spatial sampling; Prediction variance; Geostatistics

17

¹Now at: *School of Biosciences, University of Nottingham, Sutton Bonington Campus, Sutton Bonington, Leicestershire, LE12 5RD, U.K.*

**E-mail address:* murray.lark@nottingham.ac.uk (R.M. Lark).

18 **1. Introduction**

19 *1.1 The problem and its motivation*

20 How should we sample a variable in space to allow geostatistical prediction for an
21 information system or mapping project? This is an important question for the application
22 of geostatistics in soil science, particularly when limited resources are available to support
23 soil sampling in the field and the analysis of sampled material in the laboratory. It is
24 important because the sampling determines both the cost of the survey and the quality of
25 the resulting predictions.

26 One of the first approaches to this question was made by McBratney et al. (1981)
27 who showed that if the spatial covariance parameters (variogram parameters) of the target
28 random variable are known, at least approximately or from a homologous setting, then
29 one may identify the spacing of a square sample grid such that the kriging variance at
30 the centre of a grid cell (where the point kriging variance takes its largest value) does not
31 exceed some threshold. Van Groenigen et al. (1999) demonstrated that spatial simulated
32 annealing, a method for numerical optimization, can be used to find sampling designs
33 in irregularly-shaped regions so as to minimize the mean or maximum kriging variance
34 over that region. This approach will tend to produce a ‘space-filling’ or ‘spatial-coverage’
35 design, which can also be achieved by the methods of Walvoort et al (2010).

36 The limitation of spatial-coverage designs for geostatistics, be these regular grids
37 or space-filling designs in irregular regions, is that they do not provide information on
38 spatial dependence over short intervals, and so the modelled spatial covariance at short
39 lag distances is poorly constrained. The covariance at short distances is particularly
40 influential on the kriging weights. While some early geostatistical studies in soil science
41 used regular sampling grids (e.g. Burgess and Webster, 1980; Webster and Oliver, 1989)
42 it was realized that it is necessary to include some observations within a sample array that
43 are a short distance apart to support the estimation of spatial covariance parameters (e.g.
44 Atteia et al., 1994; Cattle et al., 2002). However, we are not aware of an explicit analysis

45 of the benefits of doing this in terms of the quality of final kriging predictions. Stein
46 (1999), in a simple 1-D simulation with only 20 sample locations on a regular transect,
47 showed that the likelihood function for spatial covariance parameters was very flat near
48 the maximum, but that adding just three additional observations at finer intervals within
49 the transect markedly reduced the uncertainty. In 2-D simulations with more realistic
50 sample sizes Haskard (2007) supported this finding. She considered a total sample size of
51 100, but allocated either 10 or 20 of these points to clusters within an incomplete 10×10
52 square grid. She found a marked reduction in the standard errors of spatial covariance
53 parameters when using the sample array with 10 points in a cluster by comparison to
54 the full 10×10 grid, and only a small additional benefit in using 20 of the 100 points in
55 clusters.

56 Simple spatial-coverage sampling will not do to support geostatistical prediction, so
57 how can appropriate designs be discovered? Zhu and Stein (2006) and Marchant and Lark
58 (2007a,b) showed how to define an overall objective function for the quality of a sampling
59 design, an expected mean square error of predictions, which accounts for the two sources
60 of uncertainty in the empirical best-linear unbiased prediction (E-BLUP, equivalent to the
61 kriging prediction in the general case with no covariates and the local mean assumed to
62 be stationary). These two sources are the spatial variation of the target variable and the
63 uncertainty in the maximum likelihood (ML) estimates of the spatial covariance param-
64 eters. More detail is provided in section 1.2. The key point is that we do not assume
65 that the spatial covariance parameters are known without error, but account for their
66 uncertainty, which depends in part on the sampling design. Spatial simulated annealing
67 can then be used to minimize the mean value, or the maximum value, of this objective
68 function across a study area. The resulting designs resemble a spatial-coverage sample
69 with some additional points at shorter distances.

70 These formal methods for optimization may be complex to implement. They require
71 an approximation of the spatial covariance parameters of the target variable, or a specifi-

72 cation of their joint prior distribution. In practice the scientist who is planning a survey
73 may have a more-or-less fixed sample size to deploy, and simple rules of thumb may be
74 more useful than complex procedures for optimization, which may also be computationally
75 demanding. There are various rules of thumb in geostatistics which have been influential
76 amongst practitioners. For example, it is generally advised to form empirical estimates
77 of the variogram for lag distances no longer than $D/2$ where D is the maximum distance
78 between observations (Journel and Huijbregts, 1978). Webster and Oliver (1992) suggest
79 that at least 100 observations are required to obtain a reliable estimate of the variogram.
80 Kerry et al. (2010) advise that a sampling grid for geostatistical prediction should have
81 a spacing no coarser than half the range of spatial dependence of the target variable, and
82 ideally one third to two fifths of the range.

83 The objective of this paper is to see whether it is possible to devise rules of thumb
84 to plan a geostatistical soil survey *de novo*. Following the observations of Stein (1999) and
85 Haskard (2007), and from the simulation results of Zhu and Stein (2006) and Marchant and
86 Lark (2007a,b), we propose that the rule for a geostatistical survey with N observations
87 is to withhold some number of these (a short-distance subset), distribute the remaining
88 N_{SC} according to a spatial-coverage design and then to insert each observation from the
89 short-distance subset into the resulting regular array at some fixed short distance, but
90 in a random direction, from a randomly selected site in the spatial-coverage subset. We
91 examine a quality measure for the resulting surveys, the mean square error of prediction as
92 computed by Marchant and Lark (2007a) which accounts both for the density of sampling
93 around a prediction site and the uncertainty of the spatial variance parameters. The
94 key question is whether a general recommendation can be made as to how many sample
95 sites to reserve for the short-distance subset. Our study is therefore one in the spirit of
96 ‘innovization’ (innovation by optimization), as discussed by Deb et al. (2014). The key
97 idea of innovization is that one seeks to discover rules which a practitioner can implement
98 which capture the key properties of solutions identified by formal optimization.

99 In the next section we review the calculation of the prediction error variance of the
 100 E-BLUP as proposed by Zhu and Stein (2006) and Marchant and Lark (2007a). The
 101 Methods section then sets out the sampling schemes and scenarios for which we evaluated
 102 this error variance. The scenarios correspond to random variables with a range of spatial
 103 covariance parameters. These include parameter sets selected from a Markov Chain Monte
 104 Carlo sample of parameters for the random effects in a linear mixed model for the variation
 105 of soil carbon content across a part of eastern lowland England with a range of contrasting
 106 land uses.

107 *1.2 The mean-square error of the empirical best linear unbiased prediction*

108 In this paper we consider the case of ordinary kriging, although the formulation of
 109 the problem extends to the more general best-linear unbiased prediction (BLUP) which
 110 includes universal kriging (or regression kriging in an approximately equivalent presenta-
 111 tion). The ordinary kriging prediction of a variable, Z at a location \mathbf{x}_0 , given q covariance
 112 parameters in $\boldsymbol{\theta}$ and n observations in $\mathbf{z} = (z(\mathbf{x}_1), z(\mathbf{x}_2), \dots, z(\mathbf{x}_n))^T$ can be written as

$$\widehat{Z}(\mathbf{x}_0|\boldsymbol{\theta}) = \boldsymbol{\lambda}^T \mathbf{z}, \quad (1)$$

113 where $\boldsymbol{\lambda}$ is a vector of weights. The weights are obtained from the ordinary kriging equation

$$\mathbf{L} = \mathbf{A}^{-1} \mathbf{b}, \quad (2)$$

114 where

$$\mathbf{A} = \begin{bmatrix} \mathbf{C}, & \mathbf{1}_n \\ \mathbf{1}_n^T, & 0 \end{bmatrix}$$

115 the matrix \mathbf{C} is the covariance matrix of the n observations given their locations, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$
 116 and the covariance function with parameters in $\boldsymbol{\theta}$, $C(\mathbf{x}_i - \mathbf{x}_j|\boldsymbol{\theta})$; $\mathbf{1}_n$ is a vector length n
 117 of ones,

$$\mathbf{L} = \begin{bmatrix} \boldsymbol{\lambda} \\ \psi \end{bmatrix},$$

118 where ψ is a Lagrange multiplier. If \mathbf{c} is a vector of the covariances between the target
 119 location \mathbf{x}_0 and the observations, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, then

$$\mathbf{b} = \begin{bmatrix} \mathbf{c} \\ 1 \end{bmatrix}.$$

120 In this formulation the expected square error of the prediction, the kriging variance, is

$$\sigma_{\text{OK}}^2(\mathbf{x}_0) = \mathbf{C}(\mathbf{0}|\boldsymbol{\theta}) - \mathbf{L}^T \mathbf{b}. \quad (3)$$

121 The derivations above are based on known covariance parameters, $\boldsymbol{\theta}$. In this paper
 122 we consider a frequentist framework in which $\boldsymbol{\theta}$ is treated as fixed but unknown, with the
 123 estimate $\hat{\boldsymbol{\theta}}$ obtained by maximum likelihood, see Lark (2000) for a fuller account. The
 124 estimate $\hat{\boldsymbol{\theta}}$ is ‘plugged in’ to the equations above to give the empirical BLUP (E-BLUP).
 125 Zimmerman & Cressie (1992) considered the effect of this parameter uncertainty on the
 126 kriging prediction using a Taylor series approximation. They showed that the prediction
 127 remained approximately unbiased, but an additional component of the prediction error
 128 variance should be considered. This is

$$\begin{aligned} \tau^2(\mathbf{x}_0) &= \mathbb{E} \left[\left\{ \hat{Z}(\mathbf{x}_0|\boldsymbol{\theta}) - \hat{Z}(\mathbf{x}_0|\hat{\boldsymbol{\theta}}) \right\}^2 \right] \\ &= \sum_{i=1}^q \sum_{j=1}^q \text{Cov}(\theta_i, \theta_j) \frac{\partial \hat{Z}}{\partial \theta_i} \frac{\partial \hat{Z}}{\partial \theta_j}, \end{aligned} \quad (4)$$

129 where θ_i denotes the i th parameter in $\boldsymbol{\theta}$ and $\text{Cov}(\cdot, \cdot)$ denotes the covariance of two random
 130 terms in the brackets. Zhu and Stein (2006) and Marchant and Lark (2007a) used this as
 131 a basis for a component of the expected squared prediction error. The expected value of
 132 the term due to uncertainty in the ML estimate, $\hat{\boldsymbol{\theta}}$, is

$$\mathbb{E}[\tau^2(\mathbf{x}_0)] = \sum_{i=1}^q \sum_{j=1}^q \text{Cov}(\theta_i, \theta_j) \frac{\partial \boldsymbol{\lambda}^T}{\partial \theta_i} \mathbf{C} \frac{\partial \boldsymbol{\lambda}}{\partial \theta_j}, \quad (5)$$

133 where \mathbf{C} is the covariance matrix of the n observations, given their locations. The term $\frac{\partial \boldsymbol{\lambda}}{\partial \theta_j}$
 134 is the vector of partial derivatives of the kriging weights with respect to the j th covariance
 135 parameter in $\boldsymbol{\theta}$. Marchant and Lark (2007a) provided the following equation from which
 136 these can be obtained:

$$\frac{\partial \mathbf{L}}{\partial \theta_i} = \mathbf{A}^{-1} \left(\frac{\partial \mathbf{b}}{\partial \theta_i} - \frac{\partial \mathbf{A}}{\partial \theta_i} \mathbf{A}^{-1} \mathbf{b} \right). \quad (6)$$

137 The covariance matrix for the estimated variance parameters may be approximated by the
 138 inverse of the Fisher information matrix, \mathbf{F} , so

$$\text{Cov}(\theta_i, \theta_j) \approx \mathbf{F}^{-1}(\theta_i, \theta_j), \quad (7)$$

139 where

$$\mathbf{F} = \frac{1}{2} \text{Tr} \left[\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_i} \mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_j} \right], \quad (8)$$

140 and $\text{Tr}[\cdot]$ denotes the trace of the matrix in the brackets (Kitanidis, 1987).

141 Following Zimmermann and Cressie (1992) we may obtain an overall mean square
142 error of the prediction at \mathbf{x}_0 , $\sigma_{\text{P}}^2(\mathbf{x}_0)$ as the sum of the kriging variance, Eq. (3), and the
143 expected value of $\tau^2(\mathbf{x}_0)$ given in Eq. (5):

$$\sigma_{\text{P}}^2(\mathbf{x}_0) = \sigma_{\text{OK}}^2 + \text{E}[\tau^2(\mathbf{x}_0)]. \quad (9)$$

144 It is acknowledged that this expression is an approximation, given the Taylor Series ap-
145 proximation in Eq. (4) and the comparable assumption in the use of the Fisher Information
146 matrix to obtain $\text{Cov}(\theta_i, \theta_j)$ in Eq. (7). However, Zhu and Stein (2006) suggest that this
147 approximation is reasonable, at least for comparison between sampling designs.

148 In this paper we use Eq. (9) to compute the mean squared error of E-BLUPs from
149 samples in which a specified proportion of all observations are distributed according to a
150 spatial-coverage design with each of the remaining points added to a location a short fixed
151 distance from a randomly selected point in the spatial-coverage subset. By varying the
152 total sample size, and the numbers of points in the spatial-coverage subset we were able to
153 show how the division of total sampling effort between spatial-coverage and close points
154 affects the uncertainty in the predictions, and how this differs between random variables
155 with contrasting spatial covariance parameters.

156 **3. Materials and Methods**

157 *3.1. Sampling schemes and their implementation.*

158 We start with a fixed total sample size, N . Of these N points $N_{\text{SC}} < N$ were
159 distributed according to a spatial-coverage design within a square uniform region. In the
160 initial experiment N was set to 100 and the uniform region was 256×256 units. The
161 selection of locations for the spatial-coverage points was done with the stratify procedure

162 in the `spsosa` library for the R platform (Walvoort et al., 2010; R Core Team, 2014). This
 163 procedure uses a k -means algorithm to partition a region into k units, the centroids of
 164 which constitute a spatial-coverage sample, Walvoort et al (2010) give further details. The
 165 remaining $N - N_{SC}$ points were then each allocated at random to one of the points in the
 166 spatial-coverage subsample and placed a fixed distance ($\delta = 5$ units) from the allocated
 167 point in a random direction. The fixed distance, δ , is the ‘short’ distance included in the
 168 sampling scheme to support spatial covariance modelling. We specified a fixed distance (in
 169 a random direction) for simplicity. Figure 1 shows the mean distance between a location in
 170 the region and its nearest neighbouring sample point in the spatial coverage design. The
 171 distance between a location and its nearest neighbouring sample point is the shortest lag
 172 over which the spatial covariance is required to determine the E-BLUP prediction. The
 173 value of δ was set to a short distance relative to the values in Figure 1, about one fifth the
 174 mean distance for the denser sample schemes.

175 Having generated this sample the next objective is to estimate the maximum pre-
 176 diction error variance at unsampled locations, which is equivalent to finding the kriging
 177 variance at the centre of a regular grid cell in the procedure of McBratney et al. (1981).
 178 To do this we first found the Voronoi tessellation of the spatial-coverage sample points
 179 (the short-distance subset was excluded) using the `deldir` package in R (Turner, 2015).
 180 We then found for each vertex of the set of Voronoi polygons the longest distance to a
 181 spatial-coverage sample point in one of its adjacent polygons and then found the maximum
 182 value of this distance over all vertices within the central 150×150 unit region, denoted
 183 d_{\max} . We then found a vertex, at location $\mathbf{x}_{V_{\max}}$ and a spatial-coverage point in one of
 184 the adjoining polygons $\mathbf{x}_{SC_{\max}}$ such that the vector $\mathbf{d} = \mathbf{x}_{SC_{\max}} - \mathbf{x}_{V_{\max}}$ has Euclidean
 185 norm $|\mathbf{d}| = d_{\max}$. The vertex at $\mathbf{x}_{V_{\max}}$ is necessarily a point in the domain such that
 186 none other is further from any point in the spatial-coverage set, and so is a site where the
 187 kriging variance for a prediction from points in the spatial-coverage set is large.

188 We then selected five target locations at which the prediction error variance, σ_P^2 ,

189 was evaluated. The first location was at $\mathbf{x}_{V_{\max}}$, the fifth location was at $\mathbf{x}_{V_{\max}} - \mathbf{d}(|\mathbf{d}| +$
 190 $(\delta/2))/|\mathbf{d}|$. This latter location is on a line joining $\mathbf{x}_{V_{\max}}$ to $\mathbf{x}_{SC_{\max}}$ and is distance $\delta/2$
 191 from $\mathbf{x}_{SC_{\max}}$. The remaining three points were spaced equally on the line joining the first
 192 and fifth. The contribution to the prediction error variance from the kriging variance will
 193 be largest at the target point coincident with $\mathbf{x}_{V_{\max}}$ and will be smallest at the point
 194 closest to $\mathbf{x}_{SC_{\max}}$. We computed the prediction error variance for a random variable at
 195 these five locations, given a set of spatial covariance parameters for the random variable.

196 We computed the Fisher Information Matrix, Eq. (7), for the spatial covariance pa-
 197 rameters using all N observations, but assumed that only the N_{SC} points were available for
 198 prediction. This allows us to consider the ‘worst case’ scenario, i.e. prediction in a region
 199 where only points from the spatial-coverage sample are close by. The partial derivatives
 200 with respect to the κ and ϕ parameters of the Matérn spatial covariance model (Stein,
 201 1999) were estimated numerically using the `grad` function from the `numDeriv` package for
 202 the R platform (Gilbert and Varadhan, 2015). This was done because, as discussed by
 203 Haskard et al. (2007), the analytical solutions to these derivatives are prone to rounding
 204 error. The expected value of τ^2 was then computed at each of the five target points using
 205 Eq. (5) and (6). The corresponding kriging variances were computed with Eq. (3) and
 206 then the overall prediction error variance at each site was computed with Eq. (9). The
 207 maximum prediction error variance over the five sites was then extracted.

208 The maximum prediction error variance obtained this way for a given sample size,
 209 N , with N_{SC} distributed by spatial-coverage sampling is a random variable because the
 210 spatial-coverage sample obtained may differ from one run of the `stratify` procedure to
 211 another. For this reason we repeated this procedure 50 times and calculated the mean
 212 and 95% confidence interval for each N and N_{SC} .

213 *3.2. Scenarios. 3.2.1 Contrasting random variables.* In all calculations we considered a
 214 standard random variable with a nugget variance of 0.1 and a correlated variance of 0.9.

215 We specified an isotropic Matérn spatial correlation model (Matérn, 1986; Stein, 1999)

$$\rho(\mathbf{h}) = \{2^{\kappa-1}\Gamma(\kappa)\}^{-1} \left(\frac{|\mathbf{h}|}{\phi}\right)^{\kappa} K_{\kappa}\left(\frac{|\mathbf{h}|}{\phi}\right), \quad (10)$$

216 where \mathbf{h} is the lag vector, $\Gamma(\cdot)$ is the gamma function, K_{κ} is a modified Bessel function
217 of the second kind of order κ . The two parameters are κ , a smoothness parameter, and
218 ϕ , a distance parameter. We considered cases with three values of κ : 0.5 (equivalent to
219 the exponential variogram), 0.3 (somewhat rougher than exponential) and 0.7 (somewhat
220 smoother than exponential). All three smoothness parameters were specified in combina-
221 tion of each of four effective ranges of spatial dependence at which the correlation decays
222 to approximately 0.05, equal to $\alpha_{\kappa}\phi$ where ϕ is the distance parameter of the variogram
223 and α_{κ} is equal to 2.4, 3.0 and 3.5 respectively when κ is 0.3, 0.5 and 0.7. We specified
224 values of ϕ such that the effective range was 37.5, 50, 75 or 150 units. We do not claim
225 that this range of scenarios is exhaustive. It reflects a case with a small nugget effect
226 and with the smoothness of the random process in the vicinity of that of an exponentially
227 correlated one. Within this region of the feature space of spatial covariance parameters
228 we consider the effects of variations in the effective range and the overall sample size.

229 For any given set of spatial covariance parameters we considered three sample sizes
230 with $N = 75, 100$ and 150 . In each case we considered spatial-coverage subsets with
231 $N_{\text{SC}} = N - n$ with $n = 1, 2, \dots, 10$ and then larger numbers which differed between the
232 basic sample sizes. Figure 1 shows how the mean distance from a point in the spatial-
233 coverage sample to its nearest neighbour depends on the number of points in the spatial-
234 coverage sample.

235 We kept the overall sampling density constant, irrespective of N . To achieve this
236 the dimensions of the uniform square region of interest were adjusted so that the overall
237 sampling density was equal to $100/256^2$ in all cases. This means that differences between
238 sample sizes are not confounded with effects of sample density. The effect of sample density
239 can be examined by comparing cases with fixed sample size and different effective ranges
240 of spatial dependence.

241 *3.2.2 Random variables in linear mixed models for variation of soil properties* . Here we
242 used data from a previously-published study on spatial variation of soil at within-field scale
243 in a lowland landscape in eastern England (Lark et al., 1998). The data were collected
244 from a 6-ha field in use for research on precision agriculture management of cereal crops.
245 The soil was sampled in the spring in the presence of a winter-sown barley crop but prior to
246 any fertilizer application. The soil was sampled at 100 locations which comprised a basic
247 50-m grid (24 points) with 66 additional sample points at 10-m intervals on transects
248 aligned with the 50-m grid, and 10 sample points added to allow comparisons over 5 m.
249 At each site the soil was sampled to depth 20 cm with a screw auger. In this paper we
250 use data on the organic matter content of the soil samples, which was determined by the
251 Walkley Black method (Hesse, 1971), and the nitrate content extractable by KCl (Mengel,
252 1991). Statistical analysis was conducted on the organic matter content, expressed as a
253 percent by mass of the (dry) soil, and on the natural logarithm of nitrate content (mg
254 kg^{-1} dry soil).

255 In both cases a linear mixed model was fitted to the data with a constant mean as
256 the only fixed effect. We fitted a Matérn correlation function for the correlated random
257 effect; the parameters of this function, along with a nugget variance and the variance of
258 the spatially correlated random effects, were estimated by ML.

259 We considered a situation in which a 16-ha square site, considered to be homologous
260 with the original field with respect to soil variation, is to be sampled to allow geostatistical
261 mapping of both variables. We assume that the total sample size is fixed at 100. We
262 followed the same procedures described in section 3.1 to generate realizations of space-
263 filling samples with between 99 and 50 sample points, and with the remaining sample
264 sites (1 to 50) distributed between sites of the space-filling design selected at random, and
265 placed 5 m from the associated site in a random direction. For each design we computed
266 the maximum prediction error variance in the same way described in section 3.1.

267 **4. Results**

268 Results are shown in Figures 2 to 5. In Figure 2 are shown results for all cases where
269 the effective range of the random variable of interest was 150 units. This is large relative to
270 the spacings between points in the spatial-coverage samples as shown in Figure 1. Figure
271 2a and 2b both show the maximum prediction error variance for random variables with the
272 parameter κ equal to 0.7 (somewhat smoother than an exponential random variable). The
273 different symbols correspond to the different sample sizes, and the horizontal bars show
274 the 95% confidence interval for each mean value. In Figure 2a the maximum prediction
275 error variance is plotted against the fraction of all N sample sites which are in the short-
276 distance subset (rather than the spatial-coverage subset) up to a maximum proportion of
277 0.4. In Figure 2b the same results are shown but plotted against the number of sample
278 points in the short-distance subset. Similarly, Figures 2c and 2d show maximum prediction
279 error variance plotted against, respectively, the proportion of sites in the short-distance
280 subset and the number for a random variable with $\kappa = 0.5$, and Figures 2e and 2f are
281 corresponding plots for the case with $\kappa = 0.3$. Figures 3, 4 and 5 show the corresponding
282 output for cases with the effective range of the random variable equal to 75, 50 and 37.5
283 units. The latter is of similar size to the distance to the nearest point in the spatial-
284 coverage samples.

285 In Figure 2 in all cases there is an initial reduction in the maximum prediction error
286 variance as a result of increasing the number of sample sites in the short-distance set
287 from one and a more gradual increase in the maximum prediction error variance as the
288 number of short-distance sites is increased much above 10 or so. Reducing κ (making the
289 random variable rougher) makes the response of the maximum prediction error variance
290 to the number of short-distance sample points more sensitive. For the random variable
291 with $\kappa = 0.3$ it can be seen that there is no improvement from increasing the number of
292 short-range sites above about six, but the curves are somewhat ‘flat-bottomed’, and there
293 is very little increase in the maximum prediction error variance if up to about 10 points
294 are allocated to the short-range subset. Figure 2e shows that increasing the proportion of

295 points in the short-distance set above about 0.1 increases the maximum prediction error
296 variance for all sample sizes. For the variables with κ equal to 0.5 or 0.7 (Figures 2a–2d)
297 the ‘flat-bottomed’ form of the plots is more pronounced, with very little increase in the
298 maximum prediction error variance as the size of the short-distance subset is increased.

299 Reducing the effective range of the random variable, other factors being held con-
300 stant, increases the maximum prediction error variance, as can be seen by comparing
301 Figure 2 with Figures 3–5. It also makes the increase in the maximum prediction error
302 variance as the short-distance subset is increased above 10–15 points more pronounced,
303 and the effect is very notable for the random variable with the shortest effective range
304 (Figure 5), although this also depends on the total sample size. Note that the range of
305 values on the ordinate of the plots (prediction error variance) is increased for shorter ef-
306 fective ranges, and that a wider range is used for Figure 3e and f than for the plots for
307 the smoother random variables with an effective range of 75 units.

308 In the case of the random variables with effective range of 75 or 50 units the max-
309 imum prediction error variance is markedly reduced on adding up to 2 points in the
310 short-range subset when κ is equal to 0.5 or 0.7, but for the rougher random variable
311 with $\kappa = 0.3$ further improvement is achieved by adding 5 to 7 points in the short-range
312 set. Using up to 10 sample points in the short-range set incurs a small penalty for the
313 smoothest random function ($\kappa = 0.7$) with the smallest sample size, but the increase in
314 the maximum squared prediction error is not large, and with a total sample size of 100 or
315 150 the increase is negligible for up to 20 or so points in the short-distance set.

316 Figure 5 shows results for the case where the effective range is 37.5, short relative to
317 the spacing between neighbouring sites in the spatial-coverage sample. Note that in many
318 cases the maximum squared prediction error exceed the *a priori* variance of the random
319 variable. Whilst there is a benefit from putting some points into the short-distance set the
320 increase in prediction error variance from putting too many into this set is very pronounced
321 for the smaller two sample sizes. Reduction in the spatial-coverage subset of points, with

322 the addition of extra points at a short distance, affects the uncertainty in the spatial
323 covariance parameters as well as the kriging variance component of the prediction error
324 variance when the range of spatial dependence is close to the spacing of the spatial-coverage
325 subset.

326 Comparing Figure 2e and 2f show that the absolute number of short-distance points
327 rather than the proportion of points in the subset determines the initial reduction in
328 the maximum prediction error variance (a short range set of 5-7 achieves the minimum
329 squared-prediction error regardless of the overall sample size. However, the increase in the
330 maximum prediction error variance with the proportion of sample points in the short-range
331 set is similar for all sample sizes as this proportion increases above 0.1. Examining all the
332 plots shows that setting the number of short-distance points to 10% of the total sample
333 size (shown by the vertical dotted line in the plots for the short-distance fraction) ensures
334 that sufficient short-distance points are included. For the random functions equivalent to
335 the exponential or rougher the curves are sufficiently flat that a 10% rule incurs no penalty
336 from reducing the spacing of the spatial-coverage sampling, and any such effect for the
337 smoothest random variable considered is very small.

338 Figure 6 shows the scaled variograms for organic matter content and for nitrate
339 content from the sampled field. In each case the value of the variogram is divided by the
340 sill to facilitate comparison. Note that while the values are similar over longer lags the
341 behaviour at short lags is rather different. In the case of nitrate content there is a large
342 nugget effect, but the parameter κ for the correlated random variable is 1.04, implying a
343 random process which is smoother than an exponentially correlated one. In the case of
344 organic matter content the parameter κ is 0.12, which implies a markedly rougher process.
345 The nugget effect in this latter case is zero.

346 Figure 7 shows the variogram models, scaled to an *a priori* (sill) variance of 1, for
347 organic matter and the log of nitrate content, as estimated by ML. The parameters are
348 listed in Table 1. Figure 8 shows the maximum prediction error variance for different

349 numbers of sample points out of 100 used for short-distance comparisons. The prediction
350 error variances are standardized by the *a priori* variance of the respective random effects.
351 In both cases there is a marked reduction in the error variance on the adjustment of the
352 sample design to include some observations at short distance. The minimum error variance
353 for organic matter predictions is with 8 observations used to allow comparisons over short
354 distances, and for nitrate content the minimum is with 17 such observations. In both
355 cases the rate of increase in the prediction error variance when more observations than the
356 optimal number are used for short-distance comparisons is markedly less than the rate of
357 increase as fewer such observations are included. That said, the reduction in the prediction
358 error variance for nitrate as more than 10 or so observations are used is very small. If
359 one was planning a survey to map both these variables, a design with 90 observations
360 distributed for spatial coverage and 10 included subsequently at short distances, would
361 not be markedly suboptimal for either variable. It is interesting that the 10% rule, which
362 seemed reasonably robust for the hypothetical examples with the κ parameter close to
363 0.5 and a nugget effect equal to one tenth of the sill, is also reasonable for these two soil
364 variables with rather different values of both parameters.

365 **5. Discussion and Conclusions**

366 These results show that the findings of Stein (1999) and Haskard (2007) that a
367 relatively small subset of short-distance points in a sample can markedly improve the es-
368 timation of the covariance parameters extends to the corollary that these short-distance
369 points can also improve the maximum prediction error variance, which reflects the uncer-
370 tainty in both the covariance parameters and the spatial variation between target points
371 for prediction and their neighbouring observations. Figure 1 shows that the distance be-
372 tween nearest neighbours in a spatial coverage sample increases relatively slowly as the
373 sample density is reduced over the range considered in this study. This, with the findings
374 of Stein (1999) and Haskard (2007) account for the asymmetry of the plots in Figures 2–5
375 with a reduction in the prediction error variance on the initial addition of a few close

376 points which is steeper than subsequent increases in the prediction error variance as the
377 number of points in the spatial coverage sample declines.

378 The practical conclusion is that it is important to include a short-distance subset.
379 With the larger sample size considered here, and particularly for random variables as rough
380 or rougher than the exponential, the potential cost of under-investing in short-distance
381 sampling is larger than the cost from degrading the spatial-coverage set by including
382 sample points in the short-distance subset, at least as long as the spatial-coverage set is
383 not too coarse relative to the effective range of the random variable.

384 While a very small number of sample points may markedly improve the maximum
385 prediction error variance in these examples, it must be recalled that these calculations are
386 done on the assumption of second-order stationarity. If just two or three short-distance
387 points are included then there is a risk that they will appear in atypical conditions, and this
388 could have a substantial effect on the estimated covariance parameters. For this reason the
389 inclusion of a rather larger short-distance set is good practice, and these results suggest
390 that using about 10% of the total sample effort in a short-distance subset is reasonable.

391 In this study we restrict the supplementation of the spatial coverage sample to
392 single points at a fixed distance from one of the spatial coverage set. In the sample
393 schemes presented by Marchant and Lark (2007a) that optimize a mean prediction error
394 variance the outputs resemble spatial coverage samples with either close pairs or, for
395 some sets of covariance parameters, short ‘transects’. This suggests that there might be
396 scope for further studies in the spirit of innovization to uncover rules which relate the
397 supplementation strategy to the properties of the underlying random variable. Whether
398 this provides a basis for practical rules of thumb would depend on, first, the sensitivity of
399 the optimal strategy to the (unknown) properties of the underlying random variable and,
400 second, how far the robustness of this strategy depends on the stationarity assumption.
401 Additional questions for further work would include whether it is more effective to include
402 supplementary points at a fixed distance from the spatial coverage points, or to include

403 them at random distances bounded by a maximum.

404 The expression for the prediction error variance used in this paper is for the case of
405 the E-BLUP where the mean is treated as an unknown constant, equivalent to ordinary
406 kriging. There is no reason why this should not be extended to a more general case
407 where the mean is modelled as a function of some environmental covariates. We would
408 not expect the general conclusions to differ much, because supplementary points will have
409 a negligible effect on estimation of fixed effects coefficients, but an investigation of the
410 question would be a useful further study.

411 To conclude, on the basis of these results we may recommend that, provided the
412 spatial-coverage subset of sample points is sufficiently dense to be reasonably confident
413 that the spatial dependence of the target variable is resolved, it is good practice to include
414 a short-distance subset and a relatively small investment of sample effort in such sample
415 points, which add little to the field effort required for sampling, can have a large effect
416 on the uncertainty of kriging predictions. In our hypothetical examples, where the nugget
417 effect is small to moderate (10%) and the smoothness parameter is close to that for the
418 exponential variogram (with both rougher and smoother conditions considered) a robust
419 strategy is to use a short-distance subset that corresponds to about 10% of the total
420 sample size for a range of values of the effective range of spatial dependence. We showed
421 how the optimal size of the short-distance subset could be found using the variogram for
422 two variables in a real data set on the soil. Interestingly, applying the 10% rule would be
423 a robust strategy for both these variables although their variograms are rather different
424 from those used in the hypothetical examples.

Acknowledgements

This paper is published with the permission of the Executive Director of the British Geological Survey (NERC).

References

- Atteia, O., Webster, R., Dubois, J.-P. 1994. Geostatistical analysis of soil contamination in the Swiss Jura. *Environmental Pollution* 86, 315–327.
- Burgess, T.M., Webster, R. 1980. Optimal interpolation and isarithmic mapping of soil properties. 1. The semi-variogram and punctual kriging. *Journal of Soil Science* 31, 315–331.
- Cattle, J.A., McBratney, A.B., Minasny, B. 2002. Kriging method evaluation for assessing the spatial distribution of urban soil lead concentration, *Journal of Environmental Quality* 31, 1576–1588.
- Deb, K., Bandaru, S., Greiner, D., Gaspar-Cunha, A., Tutum, C.C. 2014. An integrated approach to automated innovation for discovering useful design principles: case studies from engineering. *Applied Soft Computing*. 15, 42–56.
- Gilbert, P., Varadhan, R. 2015. numDeriv: Accurate Numerical Derivatives. R package version 2014.2-1. <http://CRAN.R-project.org/package=numDeriv>
- Haskard, K.A., 2007. An anisotropic Matérn spatial covariance model: REML estimation and properties. PhD Thesis, University of Adelaide, South Australia. Available at <http://hdl.handle.net/2440/47972>
- Haskard, K.A., Cullis, B.R., Verbyla, A.P., 2007. Anisotropic Matérn correlation and spatial prediction using REML. *Journal of Agricultural, Biological and Environmental Statistics* 12, 1–14.
- Hesse, P. R. 1971. *A Textbook of Soil Chemical Analysis*. London, John Murray.
- Journel, A.G., Huijbregts, C.J., 1978, *Mining Geostatistics*. Academic Press, London.
- Kerry, R., Oliver, M.A., Frogbrook, Z.L. 2010. Sampling in precision agriculture. In: *Geostatistical Applications for Precision Agriculture* (ed M.A. Oliver), pp. 35–63. Springer, Dordrecht.

- Kitanidis, P. K., 1987, Parametric estimation of covariances of regionalised variables. *Water Resources Bulletin*, 23, 557–567.
- Lark, R.M. 2000. Estimating variograms of soil properties by the method-of-moments and maximum likelihood; a comparison. *European Journal of Soil Science*, 51, 717–728.
- Lark, R.M., Catt, J.A., Stafford, J.V. 1998. Towards the explanation of within-field variability of yield of winter barley: soil series differences. *Journal of Agricultural Science*, 131, 409–416.
- Marchant, B.P., Lark, R.M. 2007a. Optimized sample schemes for geostatistical surveys. *Mathematical Geology*, 39, 113–134.
- Marchant, B.P., Lark, R.M. 2007b. The Matérn variogram model: implications for uncertainty propagation and sampling in geostatistical surveys. *Geoderma*, 140, 337–345.
- Matérn, B. 1986. *Spatial Variation*, Lecture Notes in Statistics, No. 36, Springer, New York.
- McBratney, A.B., Webster, R., Burgess, T.M. 1981. The design of optimal sampling schemes for local estimation and mapping of regionalised variables. I. Theory and Method. *Computers and Geosciences* 7, 331–334.
- Mengel, K., 1991. Available nitrogen in soils and its determination by the Nmin method and by electroultrafiltration (EUF). *Fertilizer Research*, 12, 37-52.
- R Core Team 2014. *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.
- Stein, M.L., 1999. *Interpolation of Spatial Data: Some Theory for Kriging*. Springer-Verlag, New York.

- Turner, R. 2015. *deldir*: Delaunay Triangulation and Dirichlet (Voronoi) Tessellation. R package version 0.1-9. <http://CRAN.R-project.org/package=deldir>
- van Groenigen, J.W., Siderius, W., Stein, A. 1999. Constrained optimisation of soil sampling for minimisation of the kriging variance. *Geoderma* 87, 239–259.
- Walvoort, D. J. J., Brus, D. J. and de Gruijter, J. J. 2010. An R package for spatial-coverage sampling and random sampling from compact geographical strata by k-means. *Computers & Geosciences*, 36, 1261–1267.
- Webster, R., Oliver, M.A. 1989. Optimal interpolation and isarithmic mapping of soil properties. VI. Disjunctive kriging and mapping the conditional probability. *Journal of Soil Science*, 40, 497–512.
- Webster, R., Oliver, M.A. 1992. Sample adequately to estimate variograms of soil properties. *Journal of Soil Science*, 43, 177–192.
- Williams, D.E. 1949. A rapid manometric method for the determination of carbonate in soil. *Soil Science Society of America Proceedings*, 25, 248–250.
- Zhu, Z., Stein, M.L. 2006. Spatial sampling design for prediction with estimated parameters. *Journal of Agricultural, Biological and Environmental Statistics*, 11, 24–44.
- Zimmerman, D.L., Cressie, N. 1992. Mean squared prediction error in the spatial linear model with estimated covariance parameters. *Annals of the Institute of Statistical Mathematics*, 44, 27–43.

Table 1. Parameters of the variograms estimated by maximum likelihood for soil organic matter content and (log-transformed) nitrate content.

Variable	ϕ /m	κ	Nugget variance	Correlated variance
Organic matter /%	97.9	0.12	0.00	0.92
Nitrate /log mg kg ⁻¹	34.9	1.04	0.14	0.15

Figure captions

Figure 1. Mean distance to nearest neighbour within a set of N_{SC} points in a spatial coverage sample in a 256×256 -unit square region.

Figure 2. Maximum prediction error variance over the standard set of five prediction locations as described in the text for random variables with an effective range of 150 units and $\kappa = 0.7$ (Fig. 2a,b), 0.5 (Fig. 2c,d) or 0.3 (Fig 2e,f). The plotted value is the mean over 50 realizations of the sampling scheme with horizontal bars showing the 95% confidence interval. The total sample size is indicated by the plotted symbol 150 (\blacktriangle), 100 (\circ) or 75 (\bullet). In the left-hand column (Fig 2a, 2c and 2e) the mean value of the maximum prediction error is plotted against the proportion of the total sample size withheld from the spatial coverage subset and added as the short-distance subset (up to a maximum proportion of 0.4). In the right-hand column (Fig. 2b, 2d and 2f) the same values are plotted against the number of points in the short-distance subset.

Figure 3. Maximum prediction error variance over the standard set of five prediction locations as described in the text for random variables with an effective range of 75 units and $\kappa = 0.7$ (Fig. 3a,b), 0.5 (Fig. 3c,d) or 0.3 (Fig 3e,f). The plotted value is the mean over 50 realizations of the sampling scheme with horizontal bars showing the 95% confidence interval. The total sample size is indicated by the plotted symbol 150 (\blacktriangle), 100 (\circ) or 75 (\bullet). In the left-hand column (Fig 3a, 3c and 3e) the mean value of the maximum prediction error is plotted against the proportion of the total sample size withheld from the spatial coverage subset and added as the short-distance subset (up to a maximum proportion of 0.4). In the right-hand column (Fig. 3b, 3d and 3f) the same values are plotted against the number of points in the short-distance subset.

Figure 4. Maximum prediction error variance over the standard set of five prediction

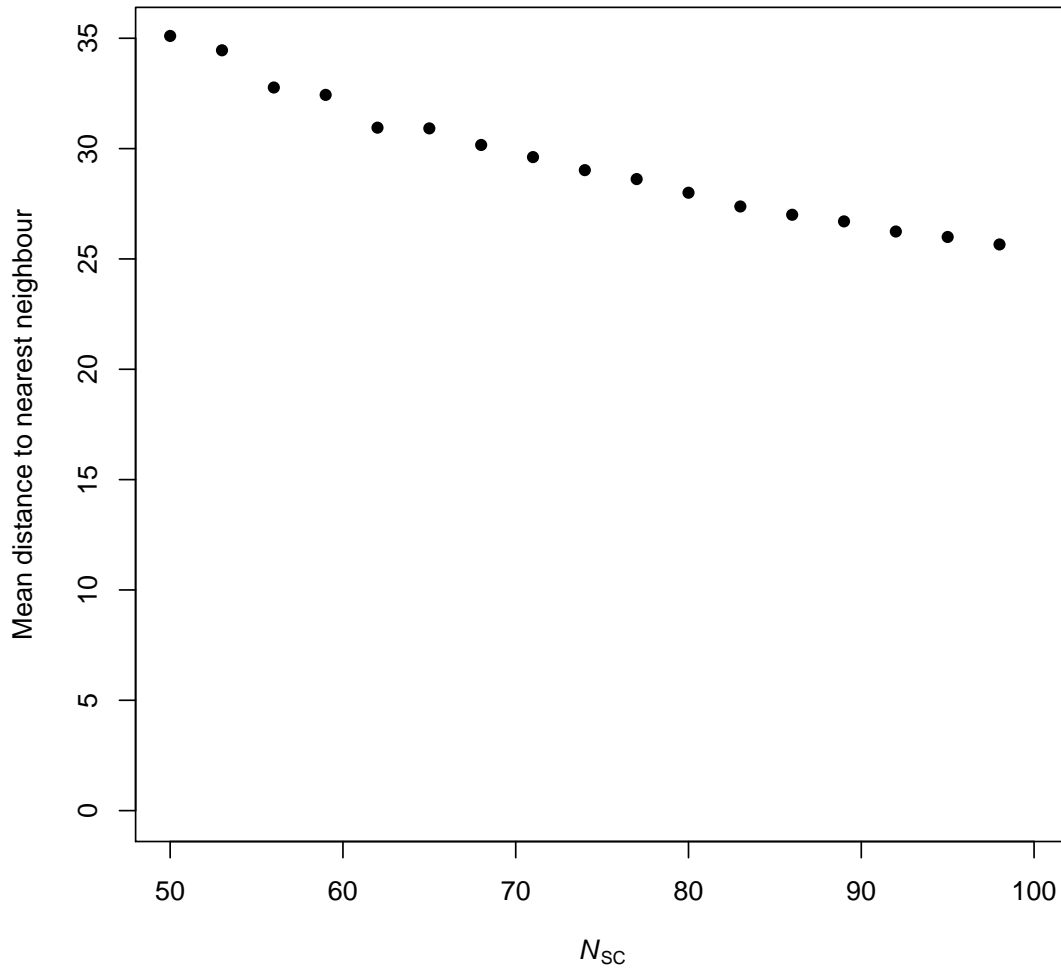
locations as described in the text for random variables with an effective range of 50 units and $\kappa = 0.7$ (Fig. 4a,b), 0.5 (Fig. 4c,d) or 0.3 (Fig 4e,f). The plotted value is the mean over 50 realizations of the sampling scheme with horizontal bars showing the 95% confidence interval. The total sample size is indicated by the plotted symbol 150 (\blacktriangle), 100 (\circ) or 75 (\bullet). In the left-hand column (Fig 4a, 4c and 4e) the mean value of the maximum prediction error is plotted against the proportion of the total sample size withheld from the spatial coverage subset and added as the short-distance subset (up to a maximum proportion of 0.4). In the right-hand column (Fig. 4b, 4d and 4f) the same values are plotted against the number of points in the short-distance subset.

Figure 5. Maximum prediction error variance over the standard set of five prediction locations as described in the text for random variables with an effective range of 37.5 units and $\kappa = 0.7$ (Fig. 5a,b), 0.5 (Fig. 5c,d) or 0.3 (Fig 5e,f). The plotted value is the mean over 50 realizations of the sampling scheme with horizontal bars showing the 95% confidence interval. The total sample size is indicated by the plotted symbol 150 (\blacktriangle), 100 (\circ) or 75 (\bullet). In the left-hand column (Fig 5a, 5c and 5e) the mean value of the maximum prediction error is plotted against the proportion of the total sample size withheld from the spatial coverage subset and added as the short-distance subset (up to a maximum proportion of 0.4). In the right-hand column (Fig. 5b, 5d and 5f) the same values are plotted against the number of points in the short-distance subset.

Figure 6. Variograms for organic matter content (broken line) and log of nitrate content (solid line) estimated by maximum likelihood and scaled to *a priori* (sill) variance of 1.

Figure 7. Mean maximum prediction error variance for organic matter content (\bullet) and log of nitrate content (\circ) with a total sample size of 100 in a 16-ha square region. The number of sample sites inserted at short distances in a space-filling design varies

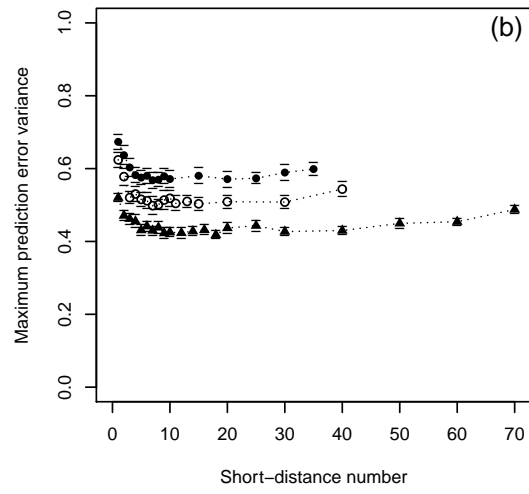
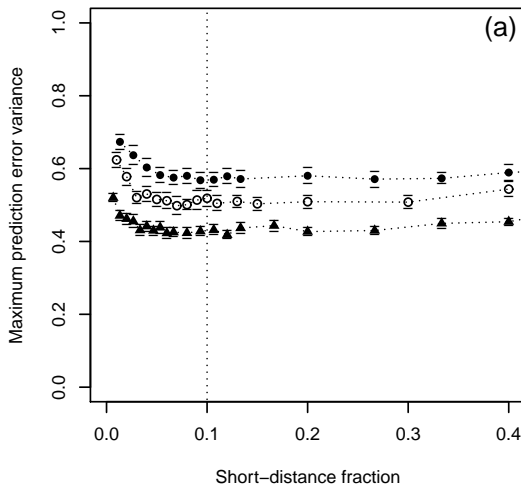
from 1 to 50. Crosses show the 95% confidence interval. The dotted vertical line shows the design where the mean maximum prediction error variance is smallest for organic matter content, and the dashed line shows the same for log-transformed nitrate content.



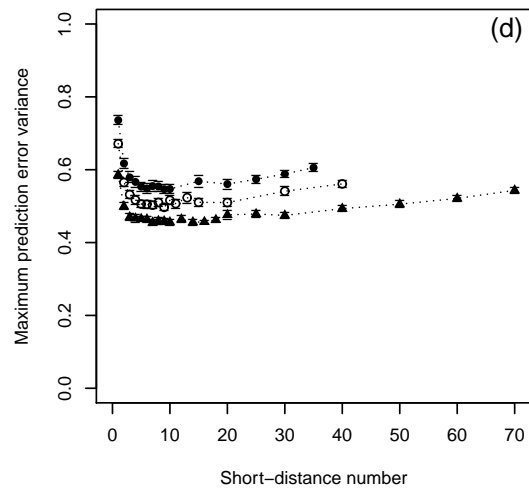
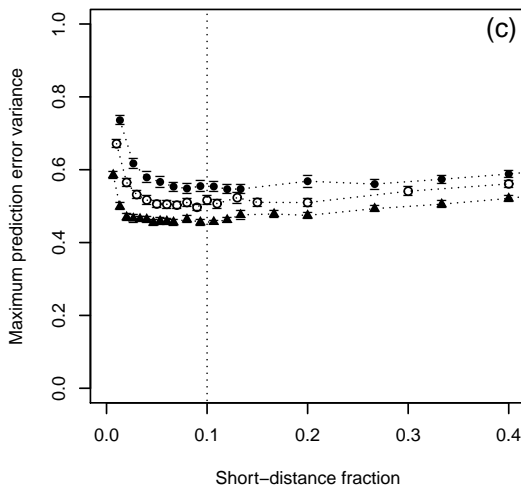
1: Fig 1

Effective range 150

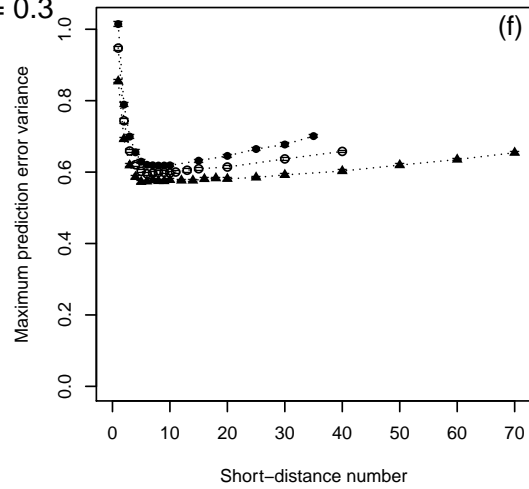
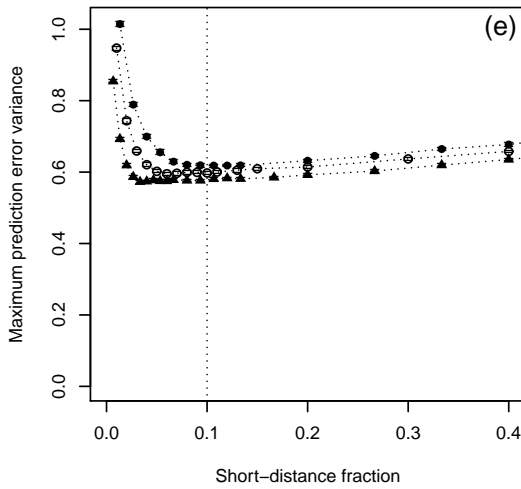
$\kappa = 0.7$



$\kappa = 0.5$

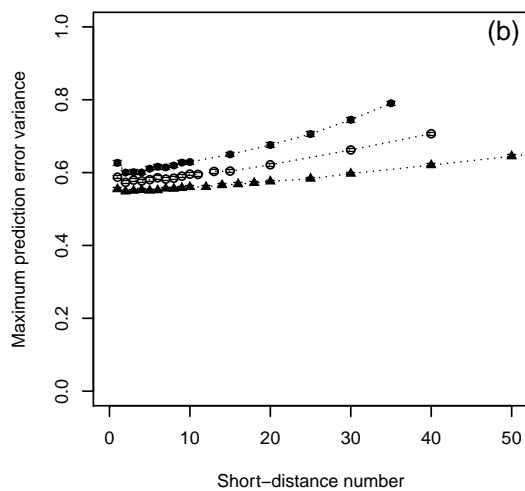
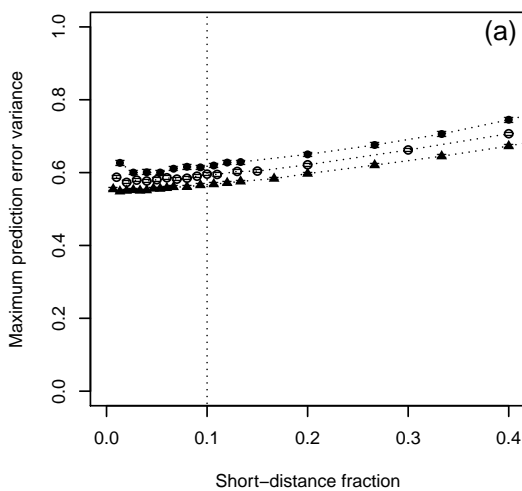


$\kappa = 0.3$

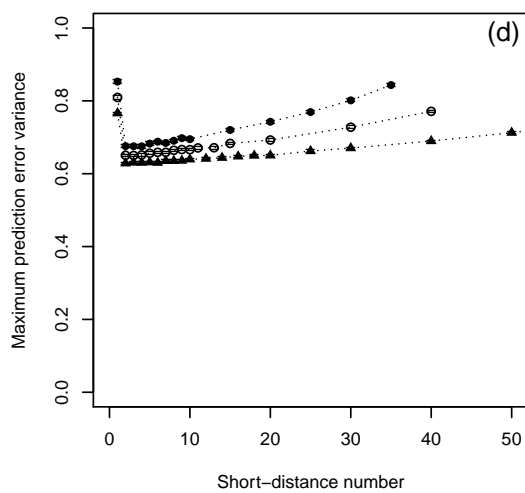
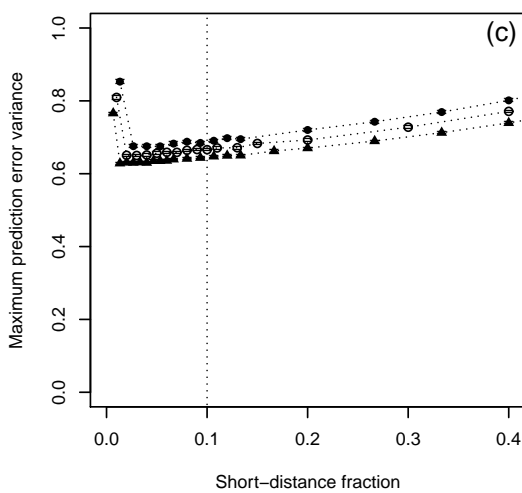


Effective range 75

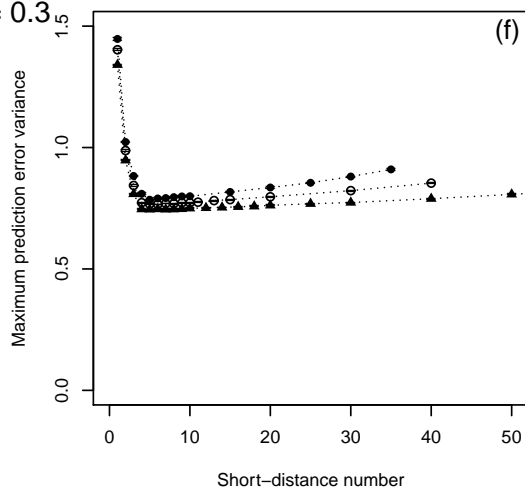
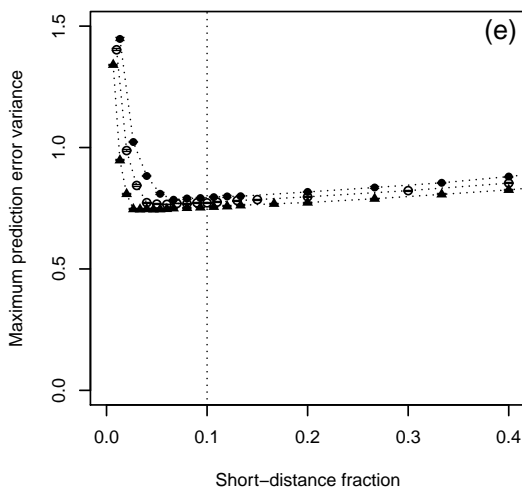
$\kappa = 0.7$



$\kappa = 0.5$

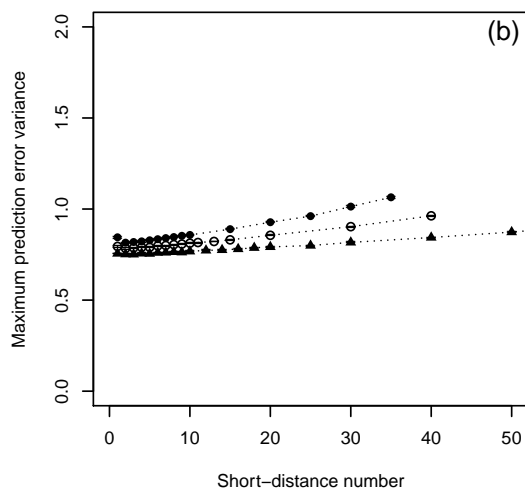
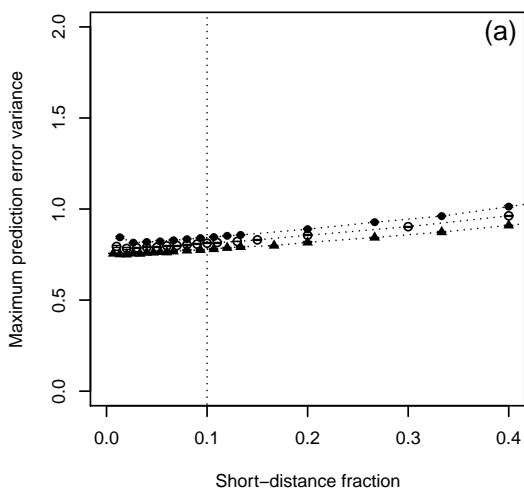


$\kappa = 0.3$

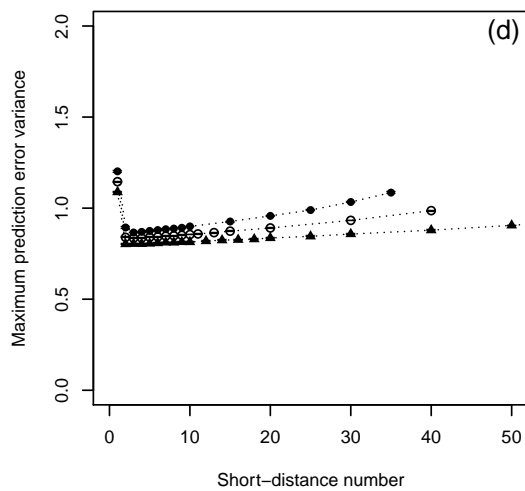
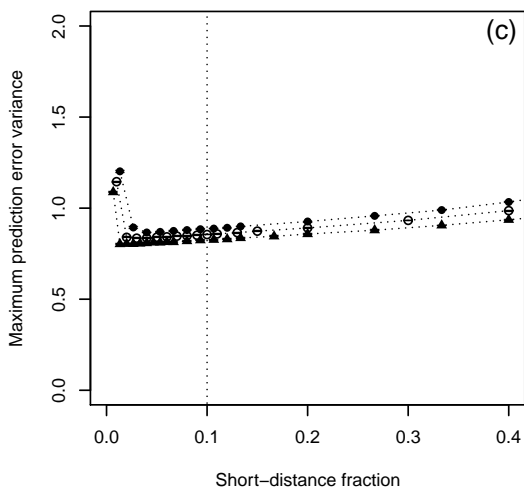


Effective range 50

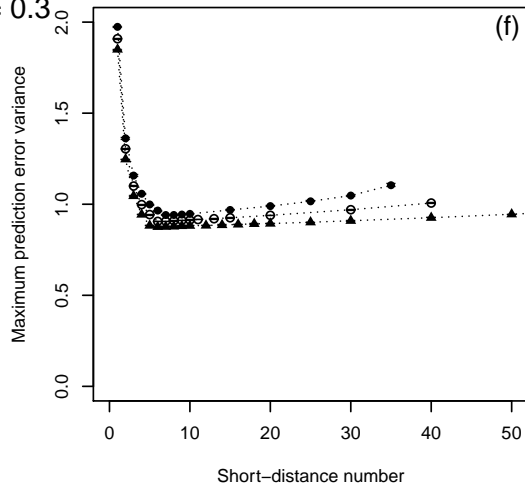
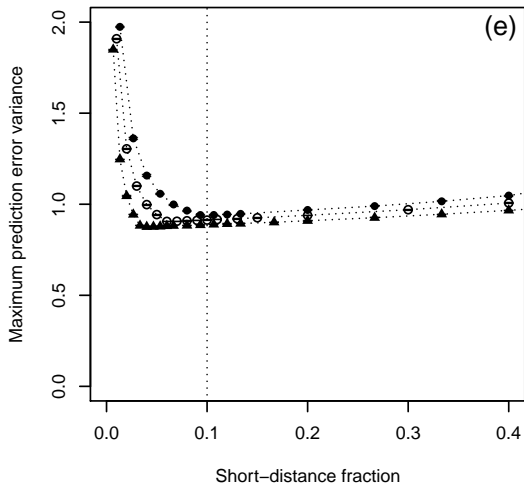
$\kappa = 0.7$



$\kappa = 0.5$

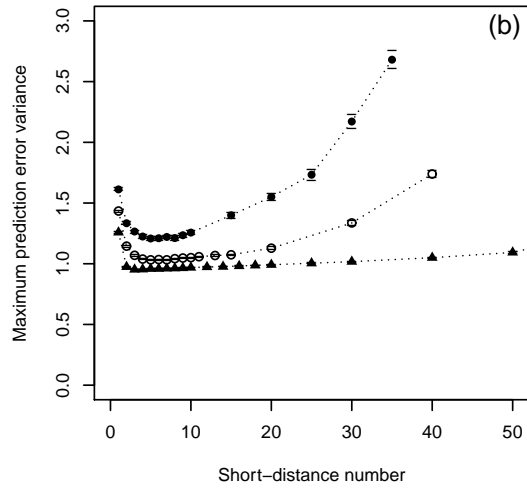
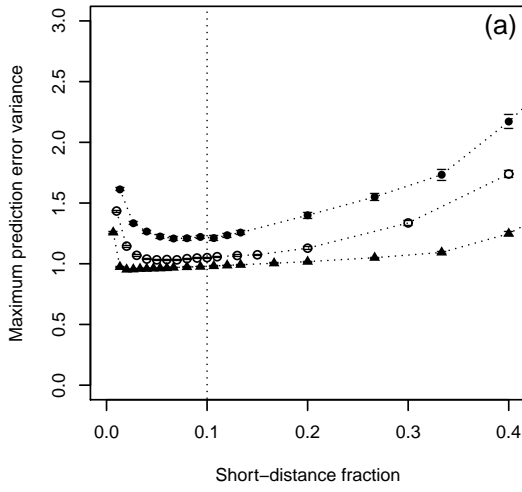


$\kappa = 0.3$

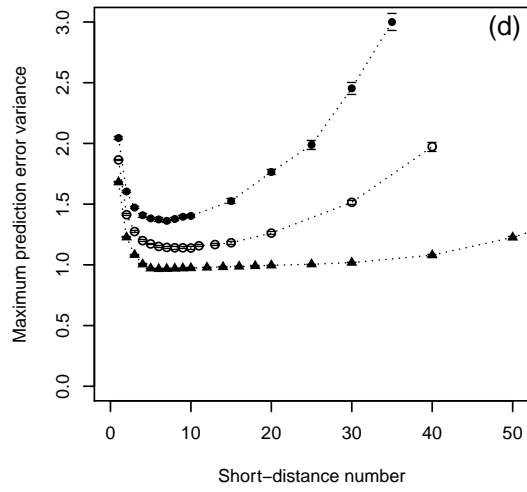
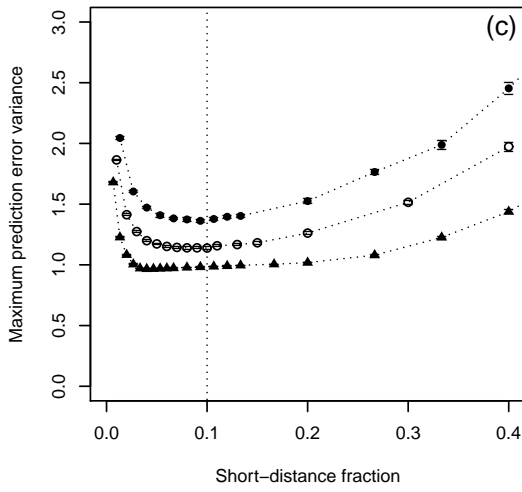


Effective range 37.5

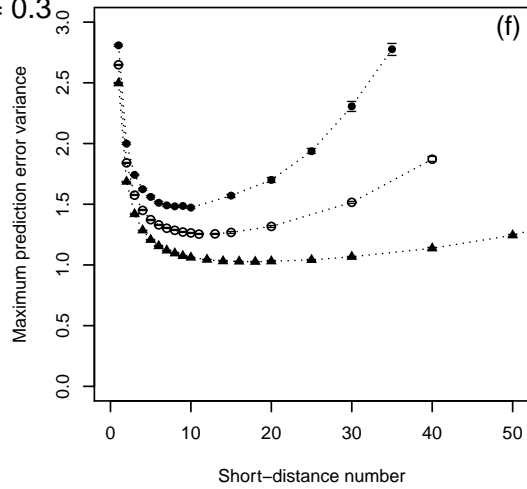
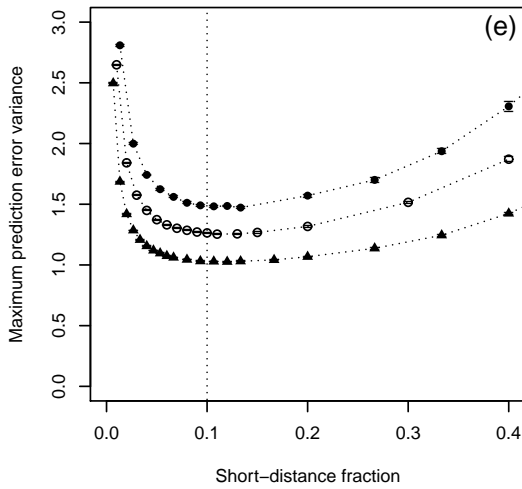
$\kappa = 0.7$

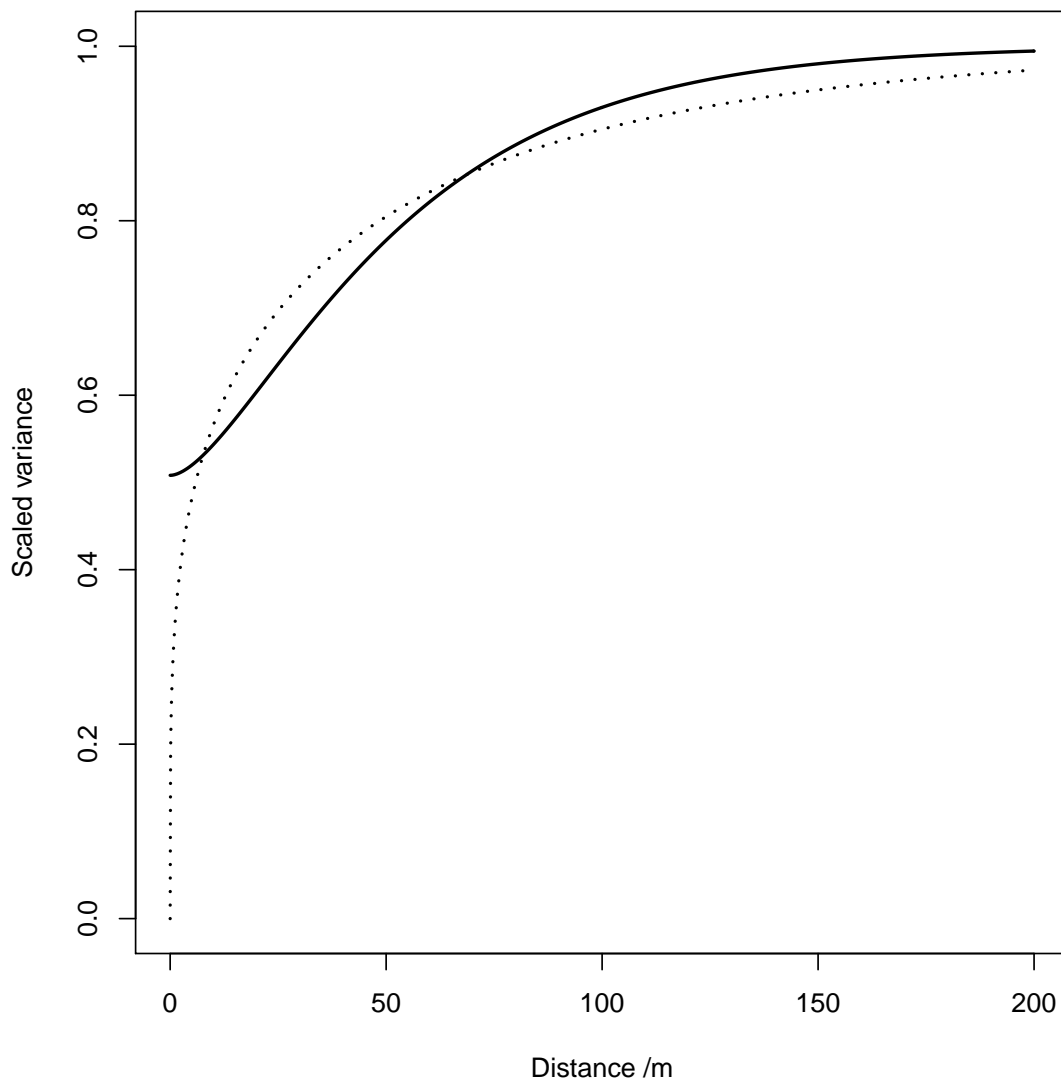


$\kappa = 0.5$

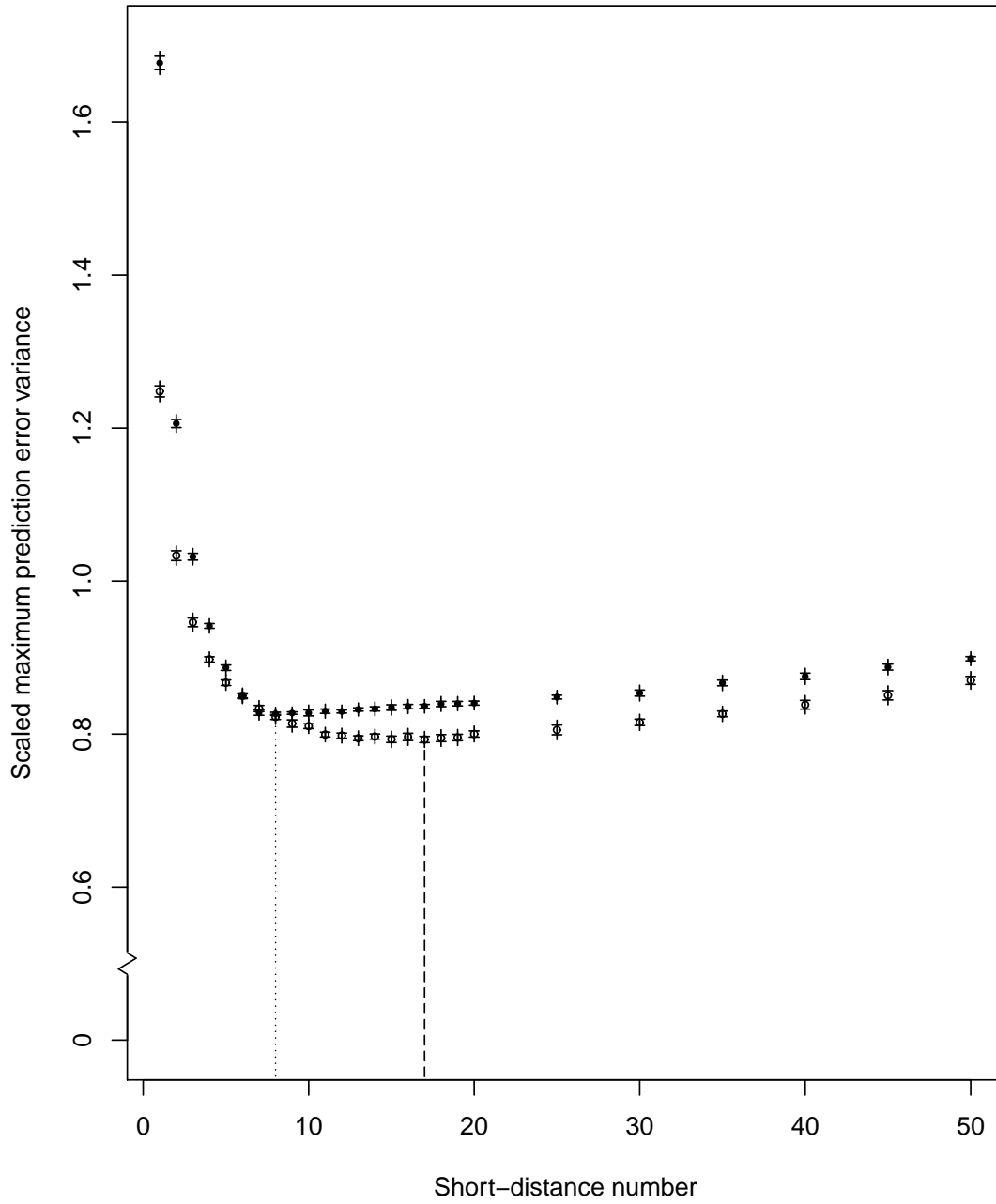


$\kappa = 0.3$





6: Fig 6



7: Fig 7